



DOCUMENTO

# IA GENERATIVA

EL CONSUMO ENERGÉTICO DE LA IA GENERATIVA

OCTUBRE DE 2024

## Contenido

El consumo energético de la IA generativa .....	2
¿Por qué la IA generativa consume tantos recursos? .....	2
Número de parámetros.....	2
Datos de entrenamiento .....	4
Complejidad de la arquitectura hardware y software .....	4
Impacto energético y medioambiental .....	5
Soluciones.....	6
Progreso algorítmico.....	6
Optimización de modelos (compresión).....	6
Nuevas arquitecturas .....	7
Fuentes de energía .....	7
Transparencia y sostenibilidad .....	8

## El consumo energético de la IA generativa

La inteligencia artificial generativa ofrece innumerables aplicaciones y usos que, hasta hace poco, parecían inalcanzables. No obstante, esta tecnología emergente enfrenta diversos desafíos y limitaciones que deberán ser superados en el futuro cercano. Uno de los principales retos es el alto consumo energético requerido, tanto en la fase de entrenamiento como durante la ejecución de los modelos.

Estudios recientes<sup>1</sup> señalan que, si la tendencia no cambia, la IA podría estar en camino de consumir anualmente tanta electricidad como todo el país de Irlanda (29,3 Tera vatios-hora por año). Este alto consumo, además, comporta la generación de CO<sub>2</sub>. Por ejemplo, el entrenamiento de Llama 3.12 en sus tres versiones (8B, 70B y 405B) ha generado 11390 toneladas de CO<sub>2</sub>.

## ¿Por qué la IA generativa consume tantos recursos?

La IA generativa consume enormes cantidades de recursos computacionales debido a varios factores, entre ellos el tamaño y complejidad de los modelos, la cantidad de datos requeridos para el entrenamiento y los cálculos involucrados en el entrenamiento y ejecución.

### Número de parámetros

El crecimiento exponencial en el número de parámetros en los modelos es uno de los principales motivos de este alto consumo de recursos. Los parámetros, como mencionamos anteriormente, son los pesos y sesgos que se ajustan durante el entrenamiento para que el modelo aprenda a generar contenido útil. Cuantos más parámetros tenga un modelo, más complejo y “capaz” es, ya que tiene una mayor capacidad para aprender patrones detallados y representaciones complejas a partir de los datos.

El número de parámetros en los modelos de IA ha crecido exponencialmente en los últimos años. Por ejemplo, el modelo GPT-2 de OpenAI (2019) tenía alrededor de 1.5 mil millones de parámetros, mientras que GPT-3 (2020) aumentó a 175 mil millones, y se estima (no se han publicado los datos) que GPT-4 tiene 1.76 billones de parámetros (billones de los españoles, no anglosajones 😊). Otros modelos similares open source, como Llama o Mistral, siguen esta tendencia.

---

<sup>1</sup> [https://www.cell.com/joule/abstract/S2542-4351\(23\)00365-3](https://www.cell.com/joule/abstract/S2542-4351(23)00365-3)

<sup>2</sup> [https://github.com/meta-llama/llama-models/blob/main/models/llama3\\_1/MODEL\\_CARD.md](https://github.com/meta-llama/llama-models/blob/main/models/llama3_1/MODEL_CARD.md)



## Datos de entrenamiento

Para que un modelo generativo funcione eficazmente, debe entrenarse con datos como texto, imágenes o vídeo. El proceso de entrenamiento implica ajustar los miles de millones de parámetros a lo largo de muchas iteraciones por la red neuronal. Esto requiere no solo de gran cantidad de procesamiento paralelo (usualmente en GPUs), sino también de almacenamiento de memoria y de ancho de banda para manejar los datos de entrenamiento y los estados intermedios del modelo. Además, las computadoras deben realizar operaciones matemáticas como multiplicaciones de matrices, de manera eficiente y rápida.

Para hacernos una idea de la cantidad de datos usados, un modelo como Llama 3 de 70 mil millones de parámetros (Llama 3 70B), liberado en abril de 2024, ha sido entrenado con 15 billones de tokens\* (billones de los españoles)<sup>4</sup>. Si lo comparamos con otros datasets similares<sup>5</sup>, esto serían aproximadamente 44 Terabytes de datos.

*\* Un **token** es una unidad básica de texto que los modelos de lenguaje utilizan para procesar datos. Puede ser una palabra completa, una parte de una palabra, o incluso un solo carácter, dependiendo de cómo el modelo divida el texto.*

## Complejidad de la arquitectura hardware y software

Por otro lado, el crecimiento exponencial en el número de parámetros, y el mecanismo de atención utilizado en los *transformers*, que hace que cada palabra o elemento en la entrada esté vinculado a cada otro (lo que significa que los cálculos aumentan cuadráticamente en relación con el tamaño de la secuencia), implica una mayor necesidad de recursos hardware (memoria RAM, almacenamiento, comunicaciones, etc.).

En un post reciente, Meta publica<sup>6</sup> que se han necesitado dos clústeres de 24.000 GPUs para entrenar Llama 3 (abril 2024), y que para finales de 2024 esperan contar con un clúster de 350.000 GPUs. Y aunque este incremento en el número de GPUs pudiera parecer exagerado, la organización sin ánimo de lucro Epoch AI, publica<sup>7</sup> que la capacidad de cómputo para entrenar modelos frontera se multiplica por 4 o 5 cada año.

---

<sup>4</sup> <https://ai.meta.com/blog/meta-llama-3/>

<sup>5</sup> <https://huggingface.co/datasets/HuggingFaceFW/fineweb>

<sup>6</sup> <https://engineering.fb.com/2024/03/12/data-center-engineering/building-metas-genai-infrastructure/>

<sup>7</sup> <https://epochai.org/blog/training-compute-of-frontier-ai-models-grows-by-4-5x-per-year>

Si revisamos el ‘model card’<sup>8</sup> (ficha del modelo), de Llama 3.1 (julio 2024) de Meta, veremos que se han requerido de 39,3 millones de horas de procesamiento, en el hardware más avanzado de Nvidia, las GPU H100-80GB.

## Impacto energético y medioambiental

La IA generativa tiene un **impacto energético** debido a los motivos expuestos en el apartado anterior. Tomemos como ejemplo los modelos de Llama 3.1 en sus diferentes versiones (8B, 70B y 405B). Para entrenar estos modelos, se utilizaron **39,3 millones de horas de GPU**. Este nivel de uso Meta lo traduce en una emisión de **11.390 toneladas de CO2**. Para poner esto en perspectiva, esa cantidad de emisiones es similar a las emisiones anuales de más de 2.400 vehículos de pasajeros o la energía utilizada por 1.000 hogares en un año. Aunque dada la reciente noticia<sup>9</sup> de The Guardian respecto a los grandes hyperscalers infra-reportando sus emisiones de CO2, esta cifra podría ser mucho más alta.

En un reciente artículo<sup>10</sup> de Epoch AI se analiza la viabilidad de escalar los modelos frontera al ritmo actual, y se proyecta que en 2030 el principal factor limitante va a ser la disponibilidad de energía para el entrenamiento. Aun suponiendo que estos futuros modelos se entrenaran en hardware 24 veces más energéticamente eficientes que el hardware actual, esta es la mejora conseguida hasta ahora, aun se requeriría 200x más energía de la usada por Llama 3.1, unos 6 GW de potencia. Por dar un orden de magnitud, España tiene instalada aproximadamente 7 GW de potencia nuclear.

Aunque hay menos información disponible, la ejecución de modelos (inferencia) también tiene un alto impacto. Algunas estimaciones<sup>11</sup> sugieren que en enero de 2023 OpenAI usaba 30.000 GPUs para manejar millones de peticiones diarias de sus usuarios, y que esas peticiones consumían alrededor de 1GWh cada día<sup>12</sup>. El equivalente a 33.000 hogares.

Algunos análisis<sup>13</sup> apuntan que en 2030 los centros de datos para IA consumirán el 4,5% de la energía global generada. El consumo de los centros de datos pasará de 49 GW en 2023 a 96 GW en 2026, de los cuales 40GW será para IA. Este mismo análisis también

---

<sup>8</sup> [https://github.com/meta-llama/llama-models/blob/main/models/llama3/MODEL\\_CARD.md](https://github.com/meta-llama/llama-models/blob/main/models/llama3/MODEL_CARD.md)

<sup>9</sup> <https://www.theguardian.com/technology/2024/sep/15/data-center-gas-emissions-tech>

<sup>10</sup> <https://epochai.org/blog/can-ai-scaling-continue-through-2030>

<sup>11</sup> <https://towardsdatascience.com/the-carbon-footprint-of-chatgpt-66932314627d>

<sup>12</sup> <https://www.washington.edu/news/2023/07/27/how-much-energy-does-chatgpt-use/>

<sup>13</sup> <https://www.semianalysis.com/p/ai-datacenter-energy-dilemma-race>

apunta que Europa estará atada de manos por la realidad geopolítica y las restricciones regulatorias estructurales sobre energía.

## Soluciones

Si bien la IA generativa tiene un potencial increíble, con lo mencionado anteriormente también se hace evidente la necesidad de buscar métodos más eficientes y sostenibles para entrenar y ejecutar estos modelos. Algunas posibles soluciones incluyen diseñar mejores algoritmos, optimizar los modelos existentes, diseñar nuevas arquitecturas hardware más eficientes, o utilizar fuentes de energía renovables o nuclear.

## Progreso algorítmico

Uno de los enfoques principales es mejorar la **eficiencia de los algoritmos** utilizados en la IA. Los avances en técnicas de entrenamiento, como la optimización del uso de datos y el ajuste más eficiente de hiperparámetros, pueden reducir el consumo energético. Por ejemplo, métodos como el **entrenamiento por lotes (batch training)** y el ajuste fino con menor cantidad de datos han demostrado ser efectivos para disminuir el número de cálculos requeridos (Eraíz-Fontanil & Schelle, 2022<sup>14</sup>). La innovación en algoritmos de aprendizaje profundo, como los **modelos de atención más eficientes**, también puede reducir los recursos computacionales necesarios.

## Optimización de modelos (compresión)

Para hacer los modelos más eficientes existen varias técnicas de compresión.

- **Quantization:** Este método consiste en reducir la precisión numérica de los pesos del modelo, cambiando de precisión de 32 bits a 16 o incluso 8 bits, lo que disminuye la carga computacional y el consumo de energía.
- **Pruning:** Involucra eliminar conexiones y nodos innecesarios en la red neuronal, manteniendo solo las partes más relevantes. Esto resulta en un modelo más pequeño y eficiente, reduciendo la cantidad de operaciones requeridas.
- **Knowledge Distillation:** Permite transferir el conocimiento de un modelo grande (modelo “profesor”) a uno más pequeño (modelo “estudiante”), manteniendo un rendimiento similar con un menor costo computacional.
- **Tensor Networks:** Estas técnicas permiten representar grandes redes neuronales de manera más compacta, disminuyendo el número de parámetros y cálculos necesarios para la inferencia y entrenamiento.

---

<sup>14</sup> Eraíz-Fontanil, F., & Arce, A. (2022). Progreso algorítmico en modelos  
<https://epochai.org/blog/algorithmic-progress-in-language-models>

## Nuevas arquitecturas

El desarrollo de hardware especializado ha sido clave para mejorar la eficiencia energética. Las **Unidades de Procesamiento Tensorial (TPUs)** de Google y las **GPUs** optimizadas para IA de NVIDIA son ejemplos de hardware diseñado específicamente para acelerar cálculos de redes neuronales con menor consumo de energía. También se están explorando arquitecturas novedosas como los **procesadores neuromórficos**, las **unidades de procesamiento fotónico**, o el uso de ordenadores cuánticos, que podrían ofrecer mejoras significativas en eficiencia.

## Fuentes de energía

Para reducir la huella de carbono generada por los centros de datos que entrenan y ejecutan modelos de IA generativa, muchas empresas tecnológicas están recurriendo a **fuentes de energía renovable** como la solar, eólica, geotérmica y la hidroeléctrica. Los centros de datos de compañías como Google y Microsoft están adoptando estas fuentes para minimizar el impacto ambiental. Por ejemplo, Google ha logrado que algunos de sus centros de datos operen con energía 100% renovable al implementar acuerdos de compra de energía renovable (PPA, power purchase agreements) o construir infraestructuras en ubicaciones con acceso a fuentes sostenibles, como es el caso reciente de la inversión de AWS en España por valor de 15.000 millones de euros.

El uso de energía renovable ayuda a reducir las emisiones de CO2 asociadas con el consumo energético de la IA. Sin embargo, también plantea desafíos. Las energías renovables, como la solar y la eólica, son **intermitentes** y dependen de las condiciones climáticas, lo que puede afectar la disponibilidad de energía. Por esta razón, las empresas a menudo deben invertir en sistemas de almacenamiento de energía, como baterías, para asegurar un suministro constante, lo cual puede aumentar los costos.

En los Estados Unidos, el crecimiento exponencial en el entrenamiento y explotación de la IA generativa también ha llevado a la industria tecnológica a buscar otras soluciones. Una de las tendencias emergentes es el uso de **energía nuclear** como fuente de energía para estos centros, en parte debido a su reciente clasificación como **energía limpia**. A pesar de los debates en torno a la energía nuclear, una de sus mayores ventajas es su **bajo impacto en las emisiones de carbono**. A lo largo de su ciclo de vida, la energía nuclear genera muy pocas emisiones de CO2, lo que la posiciona en los esfuerzos para reducir la huella de carbono de los centros de datos.

Por otro lado, y a diferencia de fuentes de energía renovable como la solar o la eólica, que son intermitentes, la energía nuclear proporciona una fuente de electricidad **constante**, lo que resulta crucial para los centros de datos que requieren operar 24/7 con gran estabilidad.

Adicionalmente, en los últimos años han surgido diferentes iniciativas del gobierno federal de Estados Unidos, y por actores privados. Por ejemplo, como se está planteando la modernización y reactivación de plantas nucleares cerradas<sup>15</sup>, aprovechando infraestructuras existentes. También han surgido nuevas tecnologías como los **reactores modulares pequeños (SMR, por sus siglas en inglés)**, reactores diseñados para ser más eficientes, seguros y escalables que las plantas nucleares tradicionales. Varias empresas tecnológicas y energéticas en EE.UU. están invirtiendo en el desarrollo de los SMR como por ejemplo Rolls-Royce<sup>16</sup>, lo que podría proporcionar una solución flexible y libre de CO2 para satisfacer las crecientes necesidades energéticas de la IA generativa.

Aunque la energía nuclear presenta muchas ventajas para los centros de datos, también plantea desafíos, como los **costos iniciales de construcción y de operación, la gestión de los residuos nucleares**, y las preocupaciones de seguridad por posibles accidentes, y por lo tanto la percepción pública negativa hacia esta fuente de energía.

## Transparencia y sostenibilidad

Dada la creciente preocupación por el impacto energético y ambiental de la IA generativa ha surgido la necesidad de una mayor **transparencia** de reporte de estos modelos. La transparencia implica que las organizaciones no solo revelen el tamaño y las capacidades de sus modelos, sino también los detalles sobre su consumo energético y las emisiones de carbono asociadas al entrenamiento y uso de estos modelos.

Para abordar este asunto, se han propuesto diferentes iniciativas y herramientas enfocadas en medir y reportar el consumo energético de los modelos de IA, como por ejemplo Carbontracker<sup>17</sup>. Esta herramienta mide las emisiones de CO2 asociadas con el entrenamiento de modelos de aprendizaje profundo. Carbontracker ayuda a los investigadores y desarrolladores a monitorear el consumo energético durante las etapas

---

<sup>15</sup> <https://www.xataka.com/energia/microsoft-reabrira-central-nuclear-que-lleva-cerrada-2019-quiere-para-alimentar-su-inteligencia-artificial>

<sup>16</sup> <https://www.rolls-royce.com/innovation/small-modular-reactors.aspx>

<sup>17</sup> <https://carbontracker.info/>

de entrenamiento, facilitando la identificación de prácticas más eficientes y el cálculo del impacto ambiental del modelo.

Por otro lado, inspirándose en las etiquetas de eficiencia energética como las que se ven en electrodomésticos o bombillas, se ha propuesto una iniciativa denominada “AI Energy Stars”<sup>18</sup> para clasificar los modelos de IA según su eficiencia energética. Esta clasificación proporciona una métrica estándar que evalúa el consumo de energía en función de factores como la cantidad de horas de GPU necesarias para el entrenamiento, la energía consumida durante la operación y las emisiones de carbono generadas.

Un sistema de clasificación energética permitiría a los investigadores y empresas comparar diferentes modelos y tomar decisiones más conscientes sobre qué modelos desarrollar o utilizar. Al etiquetar los modelos de IA con una clasificación similar a las de las “Energy Stars” en productos electrónicos, los desarrolladores podrían ser incentivados a optimizar sus modelos, reduciendo su impacto ambiental y promoviendo el desarrollo de IA más sostenible.

A nivel nacional, la Secretaría de Estado de Inteligencia Artificial (SEDIA), con la ayuda de la Asociación Española de Normalización (UNE), han puesto en marcha el grupo de trabajo CTN-UNE 71/SC 42/GT1 "Evaluación de la eficiencia energética de los sistemas de inteligencia artificial" que tiene como objetivo definir las especificaciones UNE destinadas a medir, evaluar y mejorar la eficiencia energética de los sistemas de IA, centradas en las fases de entrenamiento e inferencia.

---

<sup>18</sup> <https://www.nature.com/articles/d41586-024-02680-3>